

Protein databases

Henrik Nielsen

Protein databases, historical background

Swiss-Prot, <http://www.expasy.org/sprot/>

Established in 1986 in Switzerland

ExPASy (Expert Protein Analysis System)

Swiss Institute of Bioinformatics (SIB) and European Bioinformatics Institute (EBI)

PIR, <http://pir.georgetown.edu/>

Established in 1984

National Biomedical Research Foundation, Georgetown University, USA

In 2002 merged into:

UniProt, <http://www.uniprot.org/>

A collaboration between SIB, EBI and Georgetown University.



UniProt

UniProt Knowledgebase (UniProtKB)

UniProt Reference Clusters (UniRef)

UniProt Archive (UniParc)

UniProt Metagenomic and Environmental Sequence Database (UniMES)

UniProt Knowledgebase Release 2013_02 (06-Feb-13)

consists of:

UniProtKB/Swiss-Prot: Annotated manually (*curated*)

539,165 entries

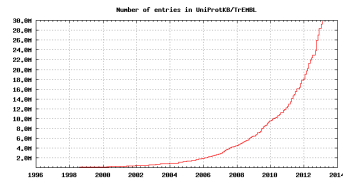
UniProtKB/TrEMBL: Computer annotated

29,769,971 entries

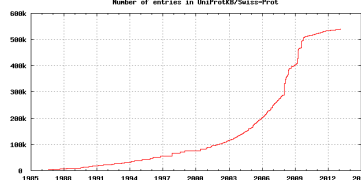
Growth of UniProt

CENTER FOR
BIOLOGICAL
CALCULATIONS
ENSCIENCE
LYSIS CBS

TrEMBL



Swiss-Prot



Content of UniProt Knowledgebase

CENTER FOR
BIOLOGICAL
CALCULATIONS
ENSCIENCE
LYSIS CBS

- Amino acid sequences
- Functional and structural annotations
 - Function / activity
 - Secondary structure
 - Subcellular location
 - Mutations, phenotypes
 - Post-translational modifications
- Origin
 - organism: Species, subspecies; classification
 - tissue
- References
- Cross references

Amino acid sequences

CENTER FOR
BIOLOGICAL
CALCULATIONS
ENSCIENCE
LYSIS CBS

From where do you get amino acid sequences?

- Translation of nucleotide sequences (GenBank/EMBL/DDBJ)
- Direct amino acid sequencing: *Edman degradation*
- Mass spectrometry
- 3D-structures

UniProt entry, formatted view

CENTER FOR
RADIOLOGICAL
CALCULATIONS
ENGENEERING
LYSIS CBS

UniProt 2.1 UniProtKB

Search

Blat * Align * Retrieve ID Mapping *

Download Contact Documentation Help

Search in

Protein Knowledgebase (UniProtKB)

Search Advanced Search + Clear

Accession #

P01009

61AT HUMAN

Reviewed UniProtKB/Swiss-Prot

Last modified: 25. Jan. 2012

Version: 101

History

Cluster with 100%, 97%, 50% identity

Documents (6)

Taskparty data

Links

Send feedback

Read comments (0) or add your own

History Annotations General annotation Ontologies Interactions All products Sequence annotation Sequences References With time

Cross refs Entry info Documents Cited by

Names and origin

Protein names

Gene names

Organism

Taxonomic identifier

Taxonomic lineage

Protein attributes

Recommended name: Alpha-1 antitrypsin

Alternative name(s): Alpha-1 protease inhibitor, Alpha-1 antiprotease, Serpin A1

Cleaved into the following chain(s): Short peptide from AAT, Short peptide from AAT

Name: SERPINA1

Synonyms: AAT, PI

ORF Names: PRO664, PRO209

Homo sapiens (human)

902 (NCBI)

Eukarya › Metazoa › Chordata › Craniata › Vertebrata › Euteleostomi › Mammalia › Eutheria › Euarchontoglires › Primates › Haplorhina › Catarrhini › Hominoidea › Homo

UniProt entry, text view (flat file)

CENTER FOR
RADIOLOGICAL
CALCULATIONS
ENGENEERING
LYSIS CBS

```
ID      A1AT_HUMAN                Reviewed:         418 AA.
AC      P01009; A6PZ14; B2RDQ8; Q0PVP3; Q13672; Q532B8; Q5U0M1; Q7M4R2;
AC      Q8GUL6; Q8GUL9; Q9G8P9; Q9G8S1; Q9P1P0; Q9UCF6; Q9UCR3;
DT      21-JUL-1986, integrated into UniProtKB/Swiss-Prot.
DT      01-OCT-1996, sequence version 3.
DT      25-JAN-2012, entry version 100.
DE      RecName: Full=Alpha-1-antitrypsin;
DE      AltName: Full=Alpha-1 protease inhibitor;
DE      AltName: Full=Alpha-1-antiproteinase;
DE      AltName: Full=Serpin A1;
DE      ContName:
DE      RecName: Full=Short peptide from AAT;
DE      Short=SPAAT;
DE      Flags: Precursor;
GN      Name=SERPINA1; Synonyms=AAT, PI; ORFNames=PRO664, PRO209;
OS      Homo sapiens (Human).
OC      Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
OC      Mammalia; Eutheria; Euarchontoglires; Primates; Haplorhini;
OC      Catarrhini; Hominoidea; Homo.
OX      NCBI TaxID=9606;
RN      [1]
RP      NUCLEOTIDE SEQUENCE [MRNA] (ISOFORM 1).
RX      MEDLINE=84107980; PubMed=6319097;
FA      Bollen A., Herzig A., Crevador A., Herion P., Chuchana P.,
FA      van der Straten A., Loriau R., Jacobs P., van Elsen A.;
RT      "Cloning and expression in Escherichia coli of full-length
RT      complementary DNA coding for human alpha 1-antitrypsin.";
RL      DNA 2:255-264 (1983).
...

```

General annotation (Comments)

CENTER FOR
RADIOLOGICAL
CALCULATIONS
ENGENEERING
LYSIS CBS

General annotation (Comments)

Function

Subcellular location

Tissue specificity

Domain

Post-translational modification

Polymorphism

Involvement in disease

Miscellaneous

Sequence similarities

Sequence caution

Inhibitor of serine proteases. Its primary target is elastase, but it also has a moderate affinity for plasmin and thrombin. Irreversibly inhibits trypsin, chymotrypsin and plasminogen activator. The aberrant form inhibits insulin-induced NO synthesis in platelets, decreases coagulation time and has proteolytic activity against insulin and plasmin. (Ref17) (Ref18) (Ref24)

Short peptide from AAT (SPAAT) is a reversible chymotrypsin inhibitor. It also inhibits elastase, but not trypsin. Its major physiological function is the protection of the lower respiratory tract against proteolytic destruction by human leukocyte elastase (HLE). (Ref17) (Ref18) (Ref24)

Secreted (Ref24)

Short peptide from AAT; Secreted › extracellular space › extracellular matrix (Ref24)

Plasma

The reactive center loop (RCL) extends out from the body of the protein and directs binding to the target protease. The protease cleaves the serpin at the reactive site within the RCL, establishing a covalent linkage between the carbonyl group of the serpin reactive site and the serine hydroxyl of the protease. The resulting inactive serpin-protease complex is highly stable.

Several isomers are observed, resulting from the combination of different N-linked glycan structures and mature N-terminus. N-linked glycan at Asn-107 is alternatively di-antennary, tri-antennary or tetra-antennary, whereas glycan at Asn-70 is di-antennary with three amounts of tri-antennary, and glycan at Asn-271 is exclusively di-antennary. The structure of the antennae is Neu5Ac(alpha1-2)Gal(beta1-4)GlcNAc attached to the core structure Man(alpha1-6)Man(alpha1-3)Man(beta1-4)GlcNAc(beta1-4)GlcNAc. Some antennae are fucosylated, which forms a Lewis-X determinant. Proteolytic processing may yield the truncated form that ranges from Asp-30 to Lys-410.

The sequence shown is that of the MTV allele which is the most common form of PI (44 to 49%). Other frequent alleles are: M1A 20 to 23%, M2 10 to 11%, M3 14 to 19%.

Defects in SERPINA1 are the cause of alpha-1-antitrypsin deficiency (A1ATD) [MIM:613400]. A disorder whose most common manifestation is emphysema, which becomes evident by the third to fourth decade. A less common manifestation of the deficiency is liver disease, which occurs in children and adults, and may result in cirrhosis and liver failure. Environmental factors, particularly cigarette smoking, greatly increase the risk of emphysema at an earlier age. (Ref20) (Ref21) (Ref22)

The aberrant form is found in the plasma of chronic smokers, and persists after smoking is ceased. It can still be found ten years after smoking has ceased.

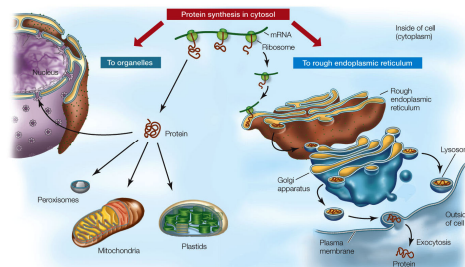
Belongs to the serpin family.

The sequence CA062334.1 differs from that shown. Reason: Erroneous initiation. Translation N-terminally shortened.

The sequence CA062396.1 differs from that shown. Reason: Erroneous initiation. Translation N-terminally shortened.

Protein sorting in eukaryotes

CENTER FOR
RIBOLOGY
CALSEQUEN
ENCEANA
LYSIS CBS



Different proteins belong to different compartments of the cell – and some belong *outside* the cell

General annotation (Comments)

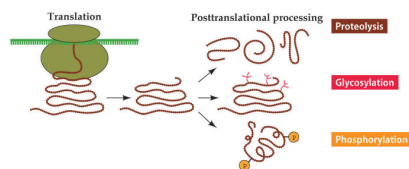
CENTER FOR
RIBOLOGY
CALSEQUEN
ENCEANA
LYSIS CBS

General annotation (Comments)

Function	Inhibitor of serine proteases. Its primary target is elastase, but it also has a moderate affinity for plasmin and thrombin. Irreversibly inhibits trypsin, chymotrypsin and plasminogen activator. The aberrant form inhibits insulin-induced NO synthesis in platelets, decreases coagulation time and has proteolytic activity against insulin and plasmin. (UniProt) (UniProt) (UniProt) Short peptide from AAT (SPAAT) is a reversible chymotrypsin inhibitor. It also inhibits elastase, but not trypsin. Its major physiological function is the protection of the lower respiratory tract against proteolytic destruction by human leukocyte elastase (HLE). (UniProt) (UniProt) (UniProt)
Subcellular location	Secreted (UniProt) Shot peptide from AAT Secreted extracellular space extracellular matrix (UniProt)
Tissue specificity	Plasma
Domain	The reactive center loop (RCL) extends out from the body of the protein and directs binding to the target protease. The protease cleaves the serpin at the reactive site within the RCL, establishing a covalent linkage between the carbonyl group of the serpin reactive site and the serine hydroxyl of the protease. The resulting inactive serpin-protease complex is highly stable.
Post-translational modification	Several isomers are observed, resulting from the combination of different N-linked glycan structure and mature N-terminus. N-linked glycan at Asn-107 is alternatively di-antennary, tri-antennary or tetra-antennary, whereas glycans at Asn-70 is di-antennary with trace amounts of tri-antennary, and glycan at Asn-271 is exclusively di-antennary. The structure of the antennae is NeuAc(alpha1-6)Gal(beta1-4)GlcNAc attached to the core structure Man(alpha1-6) (Man(alpha1-3)Man(beta1-4)GlcNAc(beta1-4)GlcNAc. Some antennae are fucosylated, which forms a Lewis-X determinant. Proteolytic processing may yield the truncated form that ranges from Asp-30 to Lys-410.
Polymorphism	The sequence shown is that of the MTV allele which is the most common form of PI (44 to 49%). Other frequent alleles are: M1A 20 to 25%, M2 10 to 11%, M3 14 to 19%.
Involvement in disease	Defects in SERPINA1 are the cause of alpha-1-antitrypsin deficiency (A1ATD) [MIM:613400]. A disorder whose most common manifestation is emphysema, which becomes evident by the third to fourth decade. A less common manifestation of the deficiency is liver disease, which occurs in children and adults, and may result in cirrhosis and liver failure. Environmental factors, particularly cigarette smoking, greatly increase the risk of emphysema at an earlier age. (UniProt) (UniProt) (UniProt)
Miscellaneous	The aberrant form is found in the plasma of chronic smokers, and persists after smoking is ceased. It can still be found ten years after smoking has ceased.
Sequence similarities	Belongs to the serpin family.
Sequence caution	The sequence CAD62334.1 differs from that shown. Reason: Erroneous initiation. Translation N-terminally shortened. The sequence CAD62595.1 differs from that shown. Reason: Erroneous initiation. Translation N-terminally shortened.

Post-translational modifications

CENTER FOR
RIBOLOGY
CALSEQUEN
ENCEANA
LYSIS CBS



Many proteins are *modified* after their synthesis in order to become active

Proteolysis: Cleavage of *signal peptides*, *propeptides* or *initiator methionine*

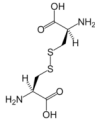
Glycosylation: Particularly common in proteins on the *surface* of cells. Also plays a role in sorting of proteins to *lysosomes*

Phosphorylation: Often *reversible*. Regulates the *activity* of many enzymes

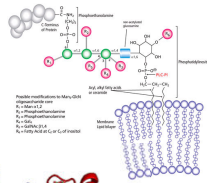
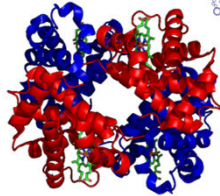
More post-translational modifications

CENTER FOR
RHEUMATOLOGY
CALIFORNIA
LYSIS CBS

- Lipid anchors
 - (e.g. GPI anchors)
- Disulfide bonds



- Prosthetic groups
 - (e.g. metal ions)



General annotation (Ontologies)

CENTER FOR
RHEUMATOLOGY
CALIFORNIA
LYSIS CBS

Ontologies

Keywords

Biological process	Acute phase Blood coagulation Extracellular matrix Secreted
Cellular component	Platelet Polysaccharide
Coding sequence diversity	Polymorphism
Domain	Signal
Molecular function	Protease inhibitor Serine protease inhibitor
PTM	Glycosylation
Technical term	Structure Complete proteome Direct protein sequencing Reference proteome

Gene Ontology (GO)

Biological process	acute phase response platelet activation platelet degranulation regulation of proteolysis extracellular space platelet alpha granule lumen proteoglycan extracellular matrix protein binding
Cellular component	extracellular space platelet alpha granule lumen proteoglycan extracellular matrix
Molecular function	protein binding



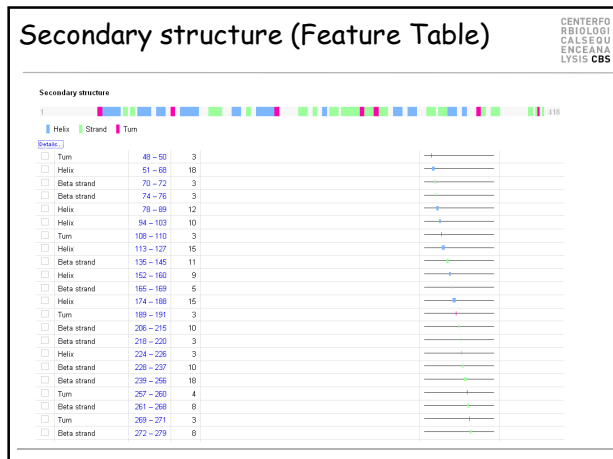
QuickGO - <http://www.ebi.ac.uk/QuickGO>

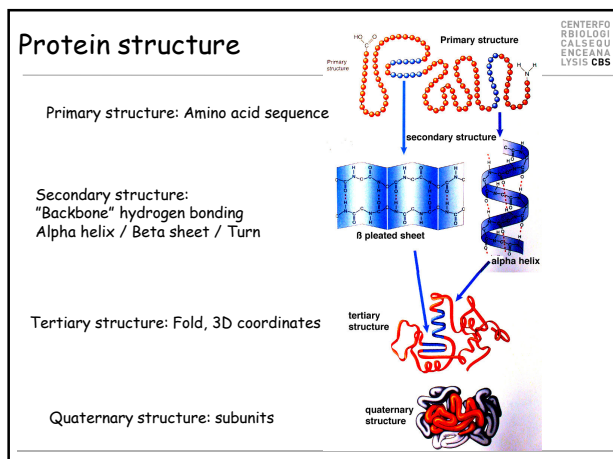
Sequence annotation (Feature Table)

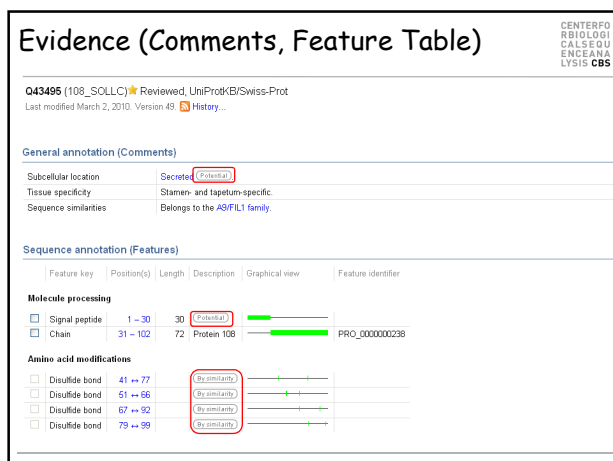
CENTER FOR
RHEUMATOLOGY
CALIFORNIA
LYSIS CBS

Sequence annotation (Features)

Feature key	Position(s)	Length	Description	Graphical view	Feature identifier
Molecule processing					
Signal peptide	1-24	24	Signal peptide (UniProt)		PRO_0000020377
Chain	25-418	394	Alpha-1-antitrypsin (UniProt)		PRO_0000034030
Peptide	375-418	44	Short peptide from AAT		
Regions					
Region	368-392	25	RCL		
Sites					
Site	362-363	2	Reactive bond		
Amino acid modifications					
Modified residue	266	1	5-cysteine (UniProt)		
Glycosylation	70	1	N-linked (GlcNAc) (UniProt)		
Glycosylation	107	1	N-linked (GlcNAc) (UniProt)		
Glycosylation	271	1	N-linked (GlcNAc) (UniProt)		
Natural variations					
Alternative sequence	307-418	112	Missing in isoform 3		VAR_000000
Alternative sequence	390-418	29	Missing in isoform 2		VAR_000000
Natural variant	4	1	S → L in ZWienham (UniProt)		VAR_000000
Natural variant	26	1	D → A in V-Munch (UniProt)		VAR_000000
Natural variant	37	1	T → A		VAR_000000
Natural variant	58	1	A → T in M-Halland (UniProt)		VAR_000000
Natural variant	63	1	R → C in (UniProt)		VAR_000000







Evidence/Confidence types

CENTER FOR
RADIOLOGICAL
CALSQUAM
ENCAANA
LYSIS CBS

3 types of *non-experimental qualifiers* in
Sequence annotation and General comment:

- *Potential*: Predicted using sequence analysis
- *Probable*: Uncertain experimental evidence
- *By similarity*: Predicted using sequence similarity

UniProt entry, sequence(s)

CENTER FOR
RADIOLOGICAL
CALSQUAM
ENCAANA
LYSIS CBS

Sequences

Sequence	Length	Mass (Da)	Tools
<input type="checkbox"/> isoform 1 [UniProt] Last modified: October 1, 1996, Version 3 Checksum: 701555F273B7F16	FASTA	418	46,737
<pre>10 20 30 40 50 60 RPSYRSGGIL LLAGLCLLYP VSLAEIPQGD AAGQCTYSHB DQDPTFNGI TYHLAEFAPS 70 80 90 100 110 120 LYDGLAQDQR ENLFFPVPY DAFAPALGSL GTDAUTDSEI LEGLNPNITE IDEALDQSP 130 140 150 160 170 180 QELLSTLRNP DQGLQLTTRH GLFLRSLGSL YQVLELVER LYHRAFPYH FQVTEAKRQ 190 200 210 220 230 240 IDPFRSGPTQ QLYGLRSLR DQUTYFALVR YIFFRQRRR FRYVDTEEE IYHNPVTYV 250 260 270 280 290 300 KYRHRSLRHR PRHSQSLGSL SHVLLNGLSL NATAITPLSD KSLQHLNDE LYHDLITDPL 310 320 330 340 350 360 EKESSRSASL BLFSLSTTGT YLAEPLQLGL GTRKYFNGA DLSQVTEAP ULKSAVHKA 370 380 390 400 VLTDSQVTE LAGNMFLEAL PHLPFVTEVP MDPVFLKEE QNTKSPFVRH KYHNPVTE</pre>			
<input type="checkbox"/> isoform 2 [UniProt] Checksum: D16A25538FB2945 Show >	FASTA	359	40,263
<input type="checkbox"/> isoform 3 [UniProt] Checksum: 15C7DBE8C25CE0C4 Show >	FASTA	306	34,756

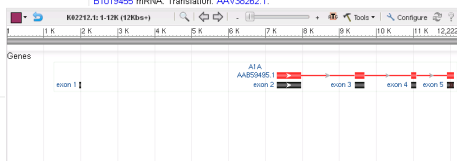
Cross-references, nucleotide sequences

CENTER FOR
RADIOLOGICAL
CALSQUAM
ENCAANA
LYSIS CBS

Sequence databases

- ☒ EMBL
- ☐ GenBank
- ☐ DDBJ

K01396 mRNA Translation: AAB59375.1.
K02217 Genomic DNA Translation: AAB59495.1.
J01853 mRNA Translation: CAA25530.1.
M11465 mRNA Translation: AAA51546.1.
J02619 Genomic DNA Translation: AAA61547.1.
D0262465 mRNA Translation: AB573300.1.
AM048535 Genomic DNA Translation: CAJ15161.1.
AF113676 mRNA Translation: AAF25581.1.
AF130065 mRNA Translation: AAG35496.1.
BX151445 mRNA Translation: CAD61914.1.
BX247988 mRNA Translation: CAD62306.1.
BX248002 mRNA Translation: CAD62334.1. Different initiation.
BX248057 mRNA Translation: CAD62555.1. Different initiation.
AK315637 mRNA Translation: BAG38005.1.
E019455 mRNA Translation: AAV38262.1.



Cross-references, 3D structure

CENTERO
RRIIOLOGI
CALISEQU
ENCEANA
LYSIS CBS

3D structure databases

- PDB
- RCSB PDB
- PDBj

Entry	Method	Resolution (Å)	Chain	Positions	PDBsum
1ATU	X-ray	2.70	A	45-418	[x]
1D55	X-ray	3.00	A	44-377	[x]
1EDX	X-ray	2.60	B	48-362	[x]
1EDX	X-ray	2.60	B	383-418	[x]
1EDX	X-ray	2.10	A	25-418	[x]
1H2Z	X-ray	2.20	A	25-418	[x]
1HCT	X-ray	3.46	A	25-418	[x]
1O08	X-ray	2.65	A	26-418	[x]
1OPH	X-ray	2.30	A	26-418	[x]
1PSI	X-ray	2.92	A	26-418	[x]
1QLP	X-ray	2.00	A	26-418	[x]
1QMB	X-ray	2.60	A	48-376	[x]
2026	X-ray	3.30	B	377-418	[x]
2026	X-ray	3.30	B	26-362	[x]
2026	X-ray	3.30	B	383-418	[x]
20UG	X-ray	2.00	A	25-418	[x]
3CWL	X-ray	2.44	A	25-418	[x]
3CWM	X-ray	2.51	A	25-418	[x]
3DRM	X-ray	2.20	A	26-418	[x]
3DRU	X-ray	3.20	A/B/C	26-418	[x]
3NDD	X-ray	1.60	A	46-372	[x]
3NDD	X-ray	1.60	B	383-418	[x]
3NDF	X-ray	2.70	A	46-381	[x]
3NDF	X-ray	2.70	B	383-418	[x]
3TIP	X-ray	3.90	A	48-418	[x]
7API	X-ray	3.00	A	36-362	[x]
8API	X-ray	3.10	B	383-418	[x]
8API	X-ray	3.10	A	36-362	[x]
9API	X-ray	3.00	B	383-418	[x]
9API	X-ray	3.00	A	36-362	[x]
9API	X-ray	3.00	B	383-418	[x]



Cross-references

CENTERO
RRIIOLOGI
CALISEQU
ENCEANA
LYSIS CBS

Other databases linked from UniProt

(there are ~100 in total):

- Nucleotide sequences
- 3D structure
- Protein-protein interactions
- Enzymatic activities and pathways
- Gene expression (microarrays and 2D-PAGE)
- Ontologies
- Families and domains
- Organism specific databases

Translation and Reading Frames

The genetic code

CENTER FOR
BIOLOGICAL
CALSEQUENCE
ANALYSIS
LYSIS CBS

		Second letter				Third letter
		U	C	A	G	
First letter	U	UUU Phenylalanine UUC Leucine UUA Leucine UUG Leucine	UCU Serine UCC Serine UCA Serine UCG Serine	UAU Tyrosine UAC Tyrosine UAA Stop codon UAG Stop codon	UGU Cysteine UGC Cysteine UGA Stop codon UGG Tryptophan	U C A G
	C	CUU Leucine CUC Leucine CUA Leucine CUG Leucine	CCU Proline CCC Proline CCA Proline CCG Proline	CAU Histidine CAC Histidine CAA Glutamine CAG Glutamine	CGU Arginine CGC Arginine CGA Arginine CGG Arginine	U C A G
	A	AUU Isoleucine AUC Isoleucine AUA Isoleucine AUG Methionine start codon	ACU Threonine ACC Threonine ACA Threonine ACG Threonine	AAU Asparagine AAC Asparagine AAA Lysine AAG Lysine	AGU Serine AGC Serine AGA Arginine AGG Arginine	U C A G
	G	GUU Valine GUC Valine GUA Valine GUG Valine	GCU Alanine GCC Alanine GCA Alanine GCG Alanine	GAU Aspartic acid GAC Aspartic acid GAA Glutamic acid GAG Glutamic acid	GGU Glycine GGC Glycine GGA Glycine GGG Glycine	U C A G

- Degenerate (*redundant*) but not ambiguous
- *Almost* universal (deviations found in mitochondria)

Reading Frames 1

CENTER FOR
BIOLOGICAL
CALSEQUENCE
ANALYSIS
LYSIS CBS

A piece of an mRNA-strand:

5' aug ccc aag cug aau agc gua gag ggg uuu uca uca uuu gag gac gau gua uaa 3'

can be divided into triplets (*codons*) in three ways:

```

1 aug ccc aag cug aau agc gua gag ggg uuu uca uca uuu gag gac gau gua uaa
  M  P  K  L  N  S  V  E  G  F  S  S  F  E  D  D  V  *
2 ugc cca agc uga aua gcg uag agg ggu uuu cau cau uug agg acg aug uau
  C  P  S  *  I  A  *  R  G  F  H  H  L  R  T  M  Y
3 gcc caa gcu gaa uag cgu aga ggg guu uuc auc auu uga gga cga ugu aua
  A  Q  A  E  *  R  R  G  V  F  I  I  *  G  R  C  I
  
```

Each possible set of triplets is called a *reading frame*.

Reading Frames 2

CENTER FOR
BIOLOGICAL
CALSEQUENCE
ANALYSIS
LYSIS CBS

Since there are two strands in DNA, there are *six* possible reading frames in a piece of DNA (three in each direction):

```

3  A  Q  A  E  *  R  R  G  V  F  I  I  *  G  R  C  I
2  C  P  S  *  I  A  *  R  G  F  H  H  L  R  T  M  Y
1  M  P  K  L  N  S  V  E  G  F  S  S  F  E  D  D  V  *
5' ATGCCCAAGCTGAATAGCTAGAGGGGTTTTCATCATTTGAGGACGATGTATAA 3'
3' TACGGGTTTCGACTTATCGCATCTCCCCAAAAGTAGTAACTCCTGCTACATATT 5'
  H  G  L  Q  I  A  Y  L  P  K  *  *  K  L  V  I  Y  L  -1
  G  L  S  F  L  T  S  P  N  E  D  N  S  S  S  T  Y  -2
  A  W  A  S  Y  R  L  P  T  K  M  M  Q  P  R  H  I  -3
  
```

A reading frame from a start codon to the first stop codon is called an *open reading frame* (underlined above).

Virtual Ribosome — required reading!

CENTER FOR
BIOLOGICAL
SEQUENCE
ANALYSIS
CBS

Nucleic Acids Research, 2006, Vol. 34, Web Server issue W385–W388
doi:10.1093/nar/gkz252

Virtual Ribosome—a comprehensive DNA translation tool with support for integration of sequence feature annotation

Rasmus Wernersson*

Center for Biological Sequence Analysis, BioCentrum-DTU, Technical University of Denmark,
Building 208, DK-2800 Lyngby, Denmark

Received February 14, 2006; Revised March 1, 2006; Accepted March 20, 2006

ABSTRACT

Virtual Ribosome is a DNA translation tool with two areas of focus. (i) Providing a strong translation tool in its own right, with an integrated ORF finder, full support for the IUPAC degenerate DNA alphabet and all translation tables defined by the NCBI taxonomy group, including the use of alternative start codons. (ii) Integration of sequence feature annotation—in particular, native support for working with files containing intron-exon structure annotation. The software is available for both download and online use at <http://www.cbs.dtu.dk/services/VirtualRibosome/>.

INTRODUCTION

A large number of software packages for translating DNA sequences already exist, as services on the World Wide

This makes it easy to build datasets that can be used for analyzing how the underlying exon structure is reflected in the protein [e.g. how exon modules map onto the 3D structure of the protein, see the ProteinMap3D server (4) elsewhere in this issue].

SOFTWARE FEATURES

Support for the degenerate nucleotide alphabet
The software has full support for the IUPAC alphabet (Table 1) for degenerate nucleotides. For example, the codon TGN correctly translates to S (serine) and not X (unknown) as often seen in other translators.

Support for a wide range of translation tables

Full support for all translation tables defined by the NCBI taxonomy group (5) (see the list below). The command-line version of the software also has support for
